

# Overlap Graph and Clumps

Mireille Régnier

LIX and INRIA

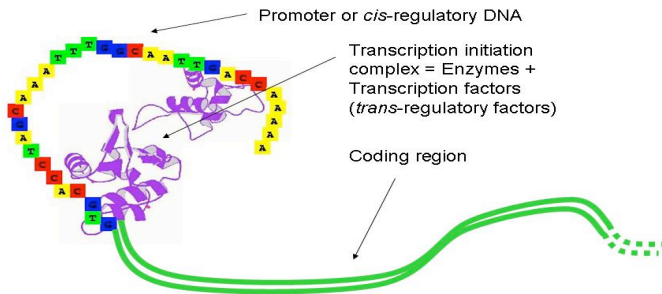
`Mireille.Regnier@inria.fr`

web page : `algo.inria.fr/regnier`

October, 9-th – 2008

- 1 Introduction and principles
- 2 Overlap Graph
- 3 Combinatorics of clumps
- 4 Open problems

The control of gene expression results from DNA/Protein interaction in the vicinity of the transcription start



# Cis-regulation changes



*Arabidopsis thaliana*



*Arabidopsis lyrata*



*Arabidopsis halleri*

# Example : the caudal motif in early developmental enhancers from *Drosophila*

GCTTTTTTATGGTCGGC  
TCGCTTTTATGGCCCAA  
CAGTTTTTATGTCTTTA  
CCGTTTTGATGGCGGTG  
AAATTTTTAGGGAACCA  
GCCCGTTTATGGTTCCC  
GACACTTTATGTGACAA  
TCGGATTTATGACACAA  
ATGTCTTTATGATTATT  
GCAACTTTTGGGCCATA  
CCCTTTTGTTGGCCCAA

Papatsenko et al., 2002

A	2	3	2	2	1	0	0	0	9	0	0	2	1	3	3	4	7
C	3	7	3	2	3	0	0	0	0	0	0	0	6	4	5	2	2
G	4	0	5	1	1	0	0	2	0	2	11	7	1	1	2	1	1
T	2	1	1	6	6	11	11	9	2	9	0	2	3	3	1	4	1

(a) Aligned Motifs

(b) Countings

# Example : the caudal motif in early developmental enhancers from *Drosophila*

GCTTTTTTATGGTCGGC  
 TCGCTTTTTATGGCCAA  
 CAGTTTTTATGTCTTTA  
 CCGTTTTGATGGCGGTG  
 AAATTTTTAGGGAACCA  
 GCCCGTTTATGGTCCC  
 GACACTTTATGTGACAA  
 TCGGATTTATGACACAA  
 ATGTCTTTATGATTATT  
 GCAACTTTTGGCCATA  
 CCCTTTTGTTGGCCAAA

Papatsenko et al., 2002

A	2	3	2	2	1	0	0	0	9	0	0	2	1	3	3	4	7
C	3	7	3	2	3	0	0	0	0	0	0	0	6	4	5	2	2
G	4	0	5	1	1	0	0	2	0	2	11	7	1	1	2	1	1
T	2	1	1	6	6	11	11	9	2	9	0	2	3	3	1	4	1

(a) Aligned Motifs

(b) Countings

A	-0.22	0.06	-0.22	-0.22	-0.62	-1.32	-1.32	-1.32	0.98	-1.32	-1.32	-0.22	-0.62	0.06	0.06	0.28	0.75
C	0.06	0.75	0.06	-0.22	0.06	-1.32	-1.32	-1.32	-1.32	-1.32	-1.32	-1.32	0.62	0.28	0.47	-0.22	-0.75
G	0.28	-1.32	0.47	-0.62	-0.62	-1.32	-1.32	-0.22	-1.32	-0.22	1.16	0.75	-0.62	-0.62	-0.22	-0.62	-0.75
T	-0.22	-0.62	-0.62	0.62	0.62	1.16	1.16	0.98	-0.22	0.98	-1.32	-0.22	0.06	0.06	6	-0.62	0.28

(c) Position Specific Scoring matrix

## Probability function

\* **Threshold  $s$**  : A word (**site**) is *similar* iff  $score(w) > s$ .

\* **Pvalue** :

$$Prob_n(\exists H; score(H) > s) .$$

## Probability function

\* **Threshold  $s$**  : A word (**site**) is *similar* iff  $score(w) > s$ .

\* **Pvalue** :

$$Prob_n(\exists H; score(H) > s) .$$

## Algorithms and data structures

\* candidates-motifs extraction



## Probability function

\* **Threshold  $s$**  : A word (**site**) is *similar* iff  $score(w) > s$ .

\* **Pvalue** :

$$Prob_n(\exists H; score(H) > s) .$$

## Algorithms and data structures

\* candidates-motifs extraction

## Model accuracy

\* Improve PWM with structural information

## Biological function

- \* Overrepresented words
- \* underrepresented words

## Statistical softwares

- \* candidates-motifs extraction
- \* statistical significance

## “Classic” methods vs Graphs

- \* induction ; [GuOd81]
- \* languages [ReSz98] ; automata [NiFISa00].

## “Classic” methods vs Graphs

- \* induction ; [GuOd81]
- \* languages [ReSz98] ; automata [NiFISa00].

## Space/time complexity

- \* Exact (all  $n$ )  $\rightarrow$  AhoPro (NII Genetika, Inria)
- \*  $O(n \times |\Sigma|)$  ;  $n$  : text size ;  $\Sigma$  : data structure.

## “Classic” methods vs Graphs

- \* induction ; [GuOd81]
- \* languages [ReSz98] ; automata [NiFISa00].

## Space/time complexity

- \* Exact (all  $n$ )  $\rightarrow$  AhoPro (NII Genetika, Inria)
- \*  $O(n \times |\Sigma|)$  ;  $n$  : text size ;  $\Sigma$  : data structure.

## Drawback

- \*  $n$  dependency ;
- \* numerical precision ;

## “Classic” methods vs Graphs

- \* induction ; [GuOd81]
- \* languages [ReSz98] ; automata [NiFISa00].

## “Classic” methods vs Graphs

- \* induction ; [GuOd81]
- \* languages [ReSz98] ; automata [NiFISa00].

## Space/time complexity

- \* Approximation  $\rightarrow$  RSA-tools, Spatt, AhoSoft (NII Genetika, Inria)
- \*  $O(1 \times |\Sigma|)$

## “Classic” methods vs Graphs

- \* induction ; [GuOd81]
- \* languages [ReSz98] ; automata [NiFISa00].

## Space/time complexity

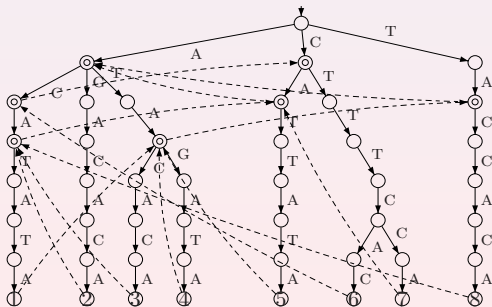
- \* Approximation  $\rightarrow$  RSA-tools, Spatt, AhoSoft (NII Genetika, Inria)
- \*  $O(1 \times |\Sigma|)$

## Drawback

- \* size of the data structure ;
- \* tightness ;

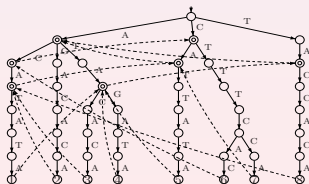


# AhoCorasick searching automaton



# AhoCorasick automaton : searching and computing

- \*  $n$  :  $w_n$  = largest prefix found = ATA;
- \*  $n + 1$  : character  $x$  found :
  - $x = G$ ,  $wx = ATAG \in Graph$ ,  $w_{n+1} = ATAG$
  - $x = A, C, T$ ,  $wx \notin Graph$ 
    - \*  $x = C$ ;  $w = A \cdot TA$ ,  $w_{n+1} = TAC \in Graph$
    - \*  $x = T$ ;  $w = AT \cdot A$ ,  $w_{n+1} = AT \in Graph$
    - \*  $x = A$ ;  $AA$ ,  $TAA \notin G$ ,  $w_{n+1} = root$



Step  $n$  :  $(p_n(w))_{w \in \text{Graph}}$ .

$p_n(w) = \text{Prob}(\text{largest prefix ending at } n \text{ is } w)$ .

## Induction

$$p_{n+1}(ATAG) = p_n(ATA) \cdot p(G)$$

$$p_{n+1}(AT) = p_n(ATA) \cdot p(T)$$

$$+ p_n(AGA) \cdot p(T)$$

$$+ p_n(CA) \cdot p(T) + p_n(TA) \cdot p(T)$$

## Left relation

$$\begin{array}{l} H_1 \mathcal{R}_L H_2 \iff \text{Father}_{\text{LOG}}(H_1) = \text{Father}_{\text{LOG}}(H_2) \\ \{ATACACA, ATAGATA\} \quad A\tilde{T}A \end{array}$$

**ATA** : Largest **prefix** of **ATACACA** that is a **suffix** in  $\mathcal{H}$

## Left relation

$$\begin{array}{l}
 H_1 \mathcal{R}_L H_2 \Leftrightarrow \text{Father}_{\text{LOG}}(H_1) = \text{Father}_{\text{LOG}}(H_2) \\
 \{ATACACA, ATAGATA\} \quad A\tilde{T}A
 \end{array}$$

**ATA** :Largest **prefix** of **ATACACA** that is a **suffix** in  $\mathcal{H}$

## Right relation

$$\begin{array}{l}
 H_1 \mathcal{R}_R H_2 \Leftrightarrow \text{Mother}_{\text{ROG}}(H_1) = \text{Mother}_{\text{ROG}}(H_2) \\
 \{ATACACA, ATACACA\} \quad A\tilde{C}A \\
 \cup \{AGACACA, \}
 \end{array}$$

**ACA** :Largest **suffix** of **ATACACA** that is a **prefix** in  $\mathcal{H}$

# Computation on Graph :induction

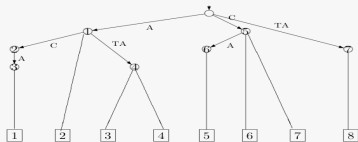


Figure 2: Left-Overlap Graph

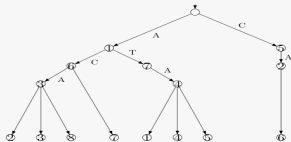


Figure 3: Right-Overlap Graph

# AhoCorasick automaton : searching and computing

First occurrence at position  $n = 18$

GGGGGGGG	ATACACA	
no $H \in \mathcal{H}$	...	$n$

# AhoCorasick automaton : searching and computing

First occurrence at position  $n = 18$

GGGGGGGG		ATACACA	
no $H \in \mathcal{H}$		...	$n$

## AND NOT

GGGG **CATT** | **ATACACA** |

GGGG **ACAT** | **ATACACA** |

GG **ACATAT** | **ATACACA** |

GG **AGACAC** | **ATACACA** |

...

All **marked** nodes in AhoGraph



Compute  $(p_n(H))_{H \in \mathcal{H}}$  using LOG, ROG.

LOG

dependency to the past

ROG

information to transfer (memory)

Compute  $(p_n(H))_{H \in \mathcal{H}}$  using LOG, ROG.

LOG

dependency to the past

ROG

information to transfer (memory)

Graph traversals...

- **First** occurrence : “small”  $n$ .
- $k$  occurrences : large  $n$ .
- $\Rightarrow$  approximation
- $\Rightarrow$  generating functions
- $\Rightarrow$  clumps

# Clump counts

With  $H_1 = \text{AACGGAA}$  and  $H_2 = \text{GAATCA}$ ,

$\text{AACGGAAACGGAACGGAATCACGGAA}$

$k$ -decomposition counted with coef.  $(-1)^k$  [BoClReVa05].

# Clump counts

With  $H_1 = \text{AACGGAA}$  and  $H_2 = \text{GAATCA}$ ,

$\text{AACGGAAACGGAACGGAATCACGGAA}$

$k$ -decomposition counted with coef.  $(-1)^k$  [BoClReVa05].

Contribution  $(-1)^7 = -1$

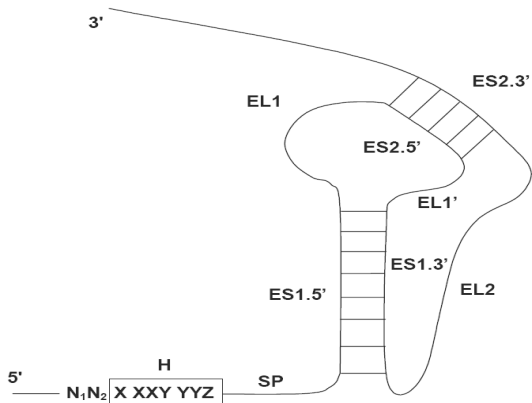
With  $\text{ACAACAACAA} = \text{AA}(\text{CAA})^3$

$\text{ACAACAACAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot$   
 $\text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{ACAACAACAA} \cdot$

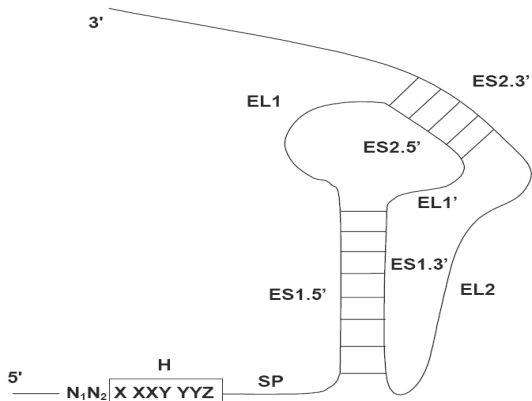
No contribution : even = odd

$\text{ACAACAACAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot$   
 $\text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{CAA} \cdot \text{ACAACAACAA} \cdot$

# Open problems : Frameshift and riboswitches



# Open problems : Frameshift and riboswitches



Boxes :

$$(w_1, w_2, \tilde{w}_1, \tilde{w}_2)$$

with : **P. Nicodeme.**

# Open problems : Frameshift and riboswitches

**A**

```

M_mazei_mtmb
M_maze_mtmb_p
M_acef_mtmb1
M_burt_mtmb1
M_burt_mtmb2
M_bark_mtmb1
M_bark_mtmb2
M_acef_mtmb2
#=GC SS_cons
UAGGGACCAGAGACUUCCCUGUCAGCUCAGGGAAAUAUAUCCUCUGACUGCAUGGGAGGCCAGAUCCAUCAGACAGCCACGAAGUUCA
UAGGGACCAGAGACUUCCCUGUCAGCUCAGGGAAAUAUAUCCUCUGACUGCAUGGGAGGCCAGAUCCAUCAGACAGCCACGAAGUUCA
UAGGGCCCCAGAGACUUCCCUGUCAGCUCAGGGAAAUAUAUCCUCUGACUGCGUGGCGGCCAGAUCCAUCAGACAGCCACGAAGUUCA
UAGGGACCAGAGACUUCCCUGUCAUCAGGGAAACAUCGCUUCGACUGUGCGGUGGCCAGAUGACAAAGUCAGCCACGAAGUUCA
UAGGGACCAGAGACUUCCCUGUCAGCUCAGGGAAAACAUCGCUUCGACUGUGCGGUGGCCAGAUGACAAAGUCAGCCACGAAGUUCA
UAGGGACCAGAGACUUCCCUGUCAGCUCAGGGAAAUAUAUCCUCUGACUGCGCGGCGGAAGUUCAAGCCAGCCAGCCACGAGGUCUG
UAGGGCCCCAGAGACUUCCCUGUCAGCUCAGGGAAAUAUAUCCUCUGACUGCGCGGCGGAAGUUCAAGCCAGCCAGCCACGAGGUCUG
UAGGGCCCCAGAGACUUCCCUGUCAGCUCAGGGAAAUAUAUCCUCUGACUGCGCGGCGGAAGUUCAAGCCAGCCAGCCACGAGGUCUG
UAGGGCCCCAGAGACUUCCCUGUCAGCUCAGGGAAAUAUAUCCUCUGACUGCGCGGCGGAAGUUCAAGCCAGCCAGCCACGAGGUCUG
UAGGGCCCCAGAGACUUCCCUGUCAGCUCAGGGAAAUAUAUCCUCUGACUGCGCGGCGGAAGUUCAAGCCAGCCAGCCACGAGGUCUG

```

